

Antoine Caubrière, Sahar Ghannay, Natalia Tomashenko, Renato De Mori, Antoine Laurent, Emmanuel Morin, Yannick Estève

ICASSP - May 2020

Error Analysis Applied to End-to-End Spoken Language Understanding

Introduction

Context

Analysis our End-to-End (E2E) Spoken Language Understanding (SLU) system This system reaches state-of-the-art performance for a french SLU task

Introduction

Context

Analysis our End-to-End (E2E) Spoken Language Understanding (SLU) system This system reaches state-of-the-art performance for a french SLU task

Goal

Analyze the errors produced by the system Understand the weakness of this E2E system From the weakness, discover how to improve our approach

Analysed system

Deep Speech 2 (DS2) [Amodei et al.] (2016)

End-to-end speech recognition system

Connectionist Temporal Classification (CTC)

Allow the system to learn the alignment between speech and output sequence to produce

Analysed system

Deep Speech 2 (DS2) [Amodei et al.] (2016)

End-to-end speech recognition system

Connectionist Temporal Classification (CTC)

Allow the system to learn the alignment between speech and output sequence to produce

End-to-End Spoken Language Understanding (SLU) [Ghannay et al.] (2018)

Tag's boundaries injection

- ASR: The sculptor Caesar died yesterday in Paris at the age of seventy-seven years
- NER : The sculptor <pers Caesar > died <time yesterday > in <loc Paris > at the age of <amount seventy-seven years >

Analysed system

Deep Speech 2 (DS2) [Amodei et al.] (2016)

End-to-end speech recognition system

Connectionist Temporal Classification (CTC)

Allow the system to learn the alignment between speech and output sequence to produce

End-to-End Spoken Language Understanding (SLU) [Ghannay et al.] (2018)

Tag's boundaries injection

- ASR: The sculptor Caesar died yesterday in Paris at the age of seventy-seven years
- NER : The sculptor <pers Caesar > died <time yesterday > in <loc Paris > at the age of <amount seventy-seven years >

Curriculum-based transfer learning (CTL) [Caubrière et al.] (2019)

Train the same model through a sequence of training processes and transfer learning

Keep all parameters except the top layer

Use of different tasks sorted from the most generic to the most specific

Automatic Speech Recognition (ASR)



Automatic Speech Recognition (ASR)

Named Entity Recognition (NER)

Annotation according to 8 entity-types (pers, loc, amount, etc)



Automatic Speech Recognition (ASR)

Named Entity Recognition (NER)

Annotation according to 8 entity-types (pers, loc, amount, etc)

Merged semantic concepts extraction (SC_mer)

MEDIA: French hotel booking task

PORTMEDIA: French theater ticket booking task

Annotation according to 76 semantic concepts (*location-town*, *stay-nbNight*, *nb-reservation*, *etc*)

Automatic Speech Recognition (ASR)

Named Entity Recognition (NER)

Annotation according to 8 entity-types (pers, loc, amount, etc)

Merged semantic concepts extraction (SC_mer)

MEDIA: French hotel booking task

PORTMEDIA: French theater ticket booking task

Annotation according to 76 semantic concepts (*location-town*, *stay-nbNight*, *nb-reservation*, *etc*)

Semantic concepts extraction on MEDIA (M)

Our target task

Automatic Speech Recognition (ASR)

Named Entity Recognition (NER)

Annotation according to 8 entity-types (pers, loc, amount, etc)

Merged semantic concepts extraction (SC_mer)

MEDIA: French hotel booking task

PORTMEDIA: French theater ticket booking task

Annotation according to 76 semantic concepts (location-town, stay-nbNight, nb-reservation, etc)

Semantic concepts extraction on MEDIA (M)

Our target task

Order of learned tasks

We define the following order of specificity: Speech > Named Entities > Semantic Concepts





French data sets

Uses of as much data as possible at our disposal Broadcast news, Telephone and Human-Human Dialogue



French data sets

Uses of as much data as possible at our disposal Broadcast news, Telephone and Human-Human Dialogue





French data sets

Uses of as much data as possible at our disposal Broadcast news, Telephone and Human-Human Dialogue

Speech ~ 360h





French data sets

Uses of as much data as possible at our disposal Broadcast news, Telephone and Human-Human Dialogue







Specific concepts





Errors distribution

Systems outputs for MEDIA concepts (development dataset)

The thirty most common mistakes



Errors distribution

Systems outputs for MEDIA concepts (development dataset)

The thirty most common mistakes



Errors distribution

Systems outputs for MEDIA concepts (development dataset)

The thirty most common mistakes



ICASSP 2020

Cases of concepts deletions (MEDIA development dataset)



Cases of concepts deletions (MEDIA development dataset)

Reference of an example

-> "From <time-date nineteen > to <time-date twenty-two > october <connectprop and > in <location-town Périgueux >"

Cases of concepts deletions (MEDIA development dataset)

Reference of an example

-> "From <time-date nineteen > to <time-date twenty-two > october <connectprop and > in <location-town Périgueux >"

Correct automatic transcription

-> "From <time-date nineteen > to <time-date twenty-two > october and in <location-town Périgueux >"

Cases of concepts deletions (MEDIA development dataset)

Reference of an example

-> "From <time-date nineteen > to <time-date twenty-two > october <connectprop and > in <location-town Périgueux >"

Correct automatic transcription

-> "From <time-date nineteen > to <time-date twenty-two > october and in <location-town Périgueux >"

Incorrect automatic transcription

-> "From <time-date nineteen > to <time-date twenty-two > october <connectprop and > par lieu "

Cases of concepts deletions (MEDIA development dataset)

Reference of an example

-> "From <time-date nineteen > to <time-date twenty-two > october <connectprop and > in <location-town Périgueux >"

Correct automatic transcription

-> "From <time-date nineteen > to <time-date twenty-two > october and in <location-town Périgueux >"

Incorrect automatic transcription

-> "From <time-date nineteen > to <time-date twenty-two > october <connectprop and > par lieu "

Correct automatic transcription but the value is nested in another concept

-> "From <time-date nineteen > to <time-date twenty-two > october <location-town and in Périgueux >"

Focused concept	Nb Deletion	Correct ASR	Wrong ASR	Nested
connectProp	39	28	6	5
lienref-coref	33	19	10	4
objet	38	31	4	3

Focused concept	Nb Deletion	Correct ASR	Wrong ASR	Nested
connectProp	39	28	6	5
lienref-coref	33	19	10	4
objet	38	31	4	3

Extra observation

Regularly ended tags without any associated started tags

-> "I'd like to know the > <object price > of the night"

A concept segmentation issue

Tackle the concept segmentation issue

Split the final MEDIA task into two tasks in the CTL approach Firstly, train the system to retrieve the boundaries of concepts only (M_seg) Secondly, train the system to specify the concepts (classical M task) CTL approach become : ASR -> NER -> SC_mer -> M_seg -> M

Tackle the concept segmentation issue

Split the final MEDIA task into two tasks in the CTL approach Firstly, train the system to retrieve the boundaries of concepts only (M_seg) Secondly, train the system to specify the concepts (classical M task) CTL approach become : ASR -> NER -> SC_mer -> M_seg -> M

Outputs to produce for M_seg

Replace each starting tag by a generic '<'

- M: "I'd like to know <lienref-coref the > <object price > of the night <connectprop and > if there are any <object rooms > left"
- M_seg: "I'd like to know < the > < price > of the night < and > if there are any < rooms > left"

Concept Error Rate (CER)

Evaluates concepts only

Concept Value Error Rate (CVER)

Evaluates concepts and values (Words within the concepts)

Concept Error Rate (CER)

Evaluates concepts only

Concept Value Error Rate (CVER)

Evaluates concepts and values (Words within the concepts)

System	CER*	CVER*
ASR -> NER -> SC_mer -> M	21.6	27.7
ASR -> NER -> SC_mer -> <i>M_seg</i> -> M	20.7	27.2
Relative gain	+4.1%	+1.0%
R : Automatic Speech Recognition	SC_mer : Merged	Semantic Concep
IER : Named Entity Recognition	M_seg : MEDI	A segmentation t
	M : MEDI	A task

Unseen Concept/Value pairs (UCV)

Examples seen in the development dataset which do not appear in the training dataset A total of 533 UCV

Unseen Concept/Value pairs (UCV)

Examples seen in the development dataset which do not appear in the training dataset A total of 533 UCV

System	Correct Concept/Value	Correct Value
ASR -> NER -> SC_mer -> M	132	38
ASR -> SC_mer -> M	124	36

Unseen Concept/Value pairs (UCV)

Examples seen in the development dataset which do not appear in the training dataset A total of 533 UCV

System	Correct Concept/Value	Correct Value
ASR -> NER -> SC_mer -> M	132	38
ASR -> SC_mer -> M	124	36

Only a small number of correct values (around a quarter) Incorrect speech transcription for a major part of the UCV

Delta of concepts recognition errors

With and without the NER task during the training



Delta of concepts recognition errors

With and without the NER task during the training



Embeddings extraction (MEDIA development dataset)

From the last bLSTM layer of DS2 DS2 trained with CTC loss function One embedding per input frame One character per input frame



Embeddings extraction (MEDIA development dataset)

From the last bLSTM layer of DS2 DS2 trained with CTC loss function One embedding per input frame One character per input frame

Words and concepts representation

Represented by more than one character Use of the sum of each frame's embeddings Use of a t-SNE transformation for a 2D representation





Observation

Each color represent a semantic class Concepts of the same classes are clustered Some very clear clusters An area with mixed concepts



Green : well-recognized concepts Red : badly-recognized concepts

Observation

Main errors are in the mixed area

Concepts errors seem to be related to an insufficiently discriminative internal representation

We presented a qualitative study of errors produced by an end-to-end SLU system

We presented a qualitative study of errors produced by an end-to-end SLU system We observed that most of the errors concern generic and domain-independent concepts

We presented a qualitative study of errors produced by an end-to-end SLU system We observed that most of the errors concern generic and domain-independent concepts We detected a concept segmentation issue for our SLU system

We presented a qualitative study of errors produced by an end-to-end SLU system We observed that most of the errors concern generic and domain-independent concepts We detected a concept segmentation issue for our SLU system We proposed an intermediate segmentation training task which allows 4.1% relative gain

We presented a qualitative study of errors produced by an end-to-end SLU system We observed that most of the errors concern generic and domain-independent concepts We detected a concept segmentation issue for our SLU system We proposed an intermediate segmentation training task which allows 4.1% relative gain We proposed a way to compute embeddings of sub-sequences

We presented a qualitative study of errors produced by an end-to-end SLU system We observed that most of the errors concern generic and domain-independent concepts We detected a concept segmentation issue for our SLU system We proposed an intermediate segmentation training task which allows 4.1% relative gain We proposed a way to compute embeddings of sub-sequences We observed that output concept errors appear to be related to an insufficiently discriminative internal representation

Perspectives



Take benefit from this cartography

How to take benefit to improve performances? How to force the system to represent the concepts in a more relevant space? Exploit the position of embeddings in the continuous space

A. Caubrière et al.

ICASSP 2020

Thank you

Contact: antoine.caubriere@univ-lemans.fr

A. Caubrière et al.

ICASSP 2020